

Report on MovieLens and Recommendation System

HarvardX Data Science Capstone Project

Raphael Kummer (GitHub)

2023-03-19

Contents

MovieLens	2
Introduction	2
Initial setup	2
Goal	2
Analysis	2
Data Inspection and preprocessing	2
Model	10
Guessing	10
Mean	10
Genre bias	10
Movie bias	10
User and Movie bias	11
User, Movie and Genre bias	11
Tuning	12
Results	13
Conclusion	13
Resources	15
System	15
Hardware	15
Software	15

MovieLens

Introduction

This is a report on the MovieLens data analysis and a recommendation model training and the models achieved performance. First the dataset is to be explored, possibly cleaned and inspected to evaluate possible training approaches. Next part is building a ML model to recommend movies to users.

Dataset

Grouplens created a movie rating dataset. The 10M dataset (Harper and Konstan 2015) used in this project is a subset of 10 million ratings of 10'000 movies by 72'000 random selected users.

Initial setup

Given is the loading of the MovieLens 10M dataset, split into an *edx* and a *final_holdout_test* set containing 10% of the MovieLens data only used for validating at the end. The dataset contains *userId*, *movieId*, *rating*, *timestamp*, *title*, and *genre*.

Goal

This dataset is used to explore and gain insight on how an effective recommendation algorithm could be developed. Such a machine learning algorithm is then developed and tested against the *final_holdout_test* set.

Analysis

Data Inspection and preprocessing

First lets take a closer look at the *edx* dataset structure and some of its content.

	userId	movieId	rating	timestamp	title	genres
1	1	122	5	838985046	Boomerang (1992)	Comedy Romance
2	1	185	5	838983525	Net, The (1995)	Action Crime Thriller
4	1	292	5	838983421	Outbreak (1995)	Action Drama Sci-Fi Thriller
5	1	316	5	838983392	Stargate (1994)	Action Adventure Sci-Fi
6	1	329	5	838983392	Star Trek: Generations (1994)	Action Adventure Drama Sci-Fi
7	1	355	5	838984474	Flintstones, The (1994)	Children Comedy Fantasy

There are several columns in this dataset:

- *userId*: an identifier for the individual user who rated the movie.
- *movieId*: an identifier for movie that was rated.
- *rating*: a rating that was given for this movie (from that user). The scale of the ratings is yet to be identified.
- *timestamp*: a timestamp in UNIX format
- *title*: the title of the movie including the year of release
- *genres*: a list of genres associated with this movie. Multiple genres are separated by '|'.

There are several aspects of the *edx* dataset to consider exploring:

- Several movies may be rated way above average because of a very good story or production.
- Some genres (or genre combinations) may be rated higher than others
- User ratings are possibly biased or generally rated higher or lower.
- User ratings in relation to a genre. User A maybe likes Horror and Action but rates movies of other genres typically lower
- User ratings in relation to movie release year

- User ratings in relation to popularity of movies (indie vs blockbuster).

And probably lots more.

The dataset *edx* contains 0 missing values.

Lets list all genres. We notice a unusual genre named (no genres listed). On further investigation, only one movie (Pull My Daisy) has no genre listed. According to IMDB the genre of this movie from 1958 is “Short”. Let’s separate the genres listed and add the first three of each movie into separate columns.

Comedy, Romance, Action, Crime, Thriller, Drama, Sci-Fi, Adventure, Children, Fantasy, War, Animation, Musical, Western, Mystery, Film-Noir, Horror, Documentary, IMAX, (no genres listed)

We have a genre with (*no genres listed*), lets inspect that:

	userId	movieId	rating	timestamp	title	genres
1025055	7701	8606	5.0	1190806786	Pull My Daisy (1958)	(no genres listed)
1453345	10680	8606	4.5	1171170472	Pull My Daisy (1958)	(no genres listed)
4066835	29097	8606	2.0	1089648625	Pull My Daisy (1958)	(no genres listed)
6456906	46142	8606	3.5	1226518191	Pull My Daisy (1958)	(no genres listed)
8046611	57696	8606	4.5	1230588636	Pull My Daisy (1958)	(no genres listed)
8988750	64411	8606	3.5	1096732843	Pull My Daisy (1958)	(no genres listed)
9404670	67385	8606	2.5	1188277325	Pull My Daisy (1958)	(no genres listed)

“Pull My Daisy” from 1958 is the only movie without a genre. According to IMDB its genre is “Short”.

As we have seen, there are 20 unique genres in the dataset.

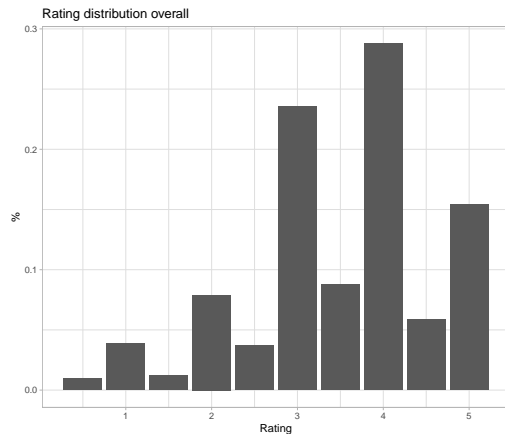
Then we transform the UNIX timestamps to date/time and for ease of use also extract the year the movie was rated. And separate the release year, embedded in the title, for possible further investigation.

Now some columns are added for further inspection, lets see the table again.

	userId	movieId	rating	timestamp	title	genres	main_genre	side1_genre	side2_genre	date	yearrated	releaseYear
1	1	122	5	838985046	Boomerang (1992)	Comedy/Romance	Comedy	Romance	NA	1996-08-02 11:24:06	1996	1992
2	1	185	5	838983525	Net, The (1995)	Action/Crime/Thriller	Action	Crime	Thriller	1996-08-02 10:58:45	1996	1995
4	1	292	5	838983421	Outbreak (1995)	Action/Drama/Sci-Fi/Thriller	Action	Drama	Sci-Fi	1996-08-02 10:57:01	1996	1995
5	1	316	5	838983392	Stargate (1994)	Action/Adventure/Sci-Fi	Action	Adventure	Sci-Fi	1996-08-02 10:56:32	1996	1994
6	1	329	5	838983392	Star Trek: Generations (1994)	Action/Adventure/Drama/Sci-Fi	Action	Adventure	Drama	1996-08-02 10:56:32	1996	1994
7	1	355	5	838984474	Flintstones, The (1994)	Children/Comedy/Fantasy	Children	Comedy	Fantasy	1996-08-02 11:14:34	1996	1994

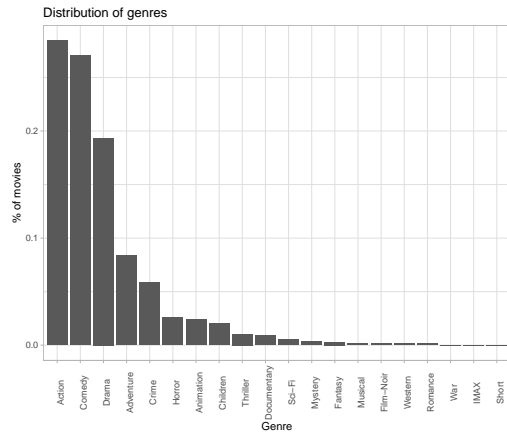
Distributions

Lets plot some distributions, starting the distribution of ratings.

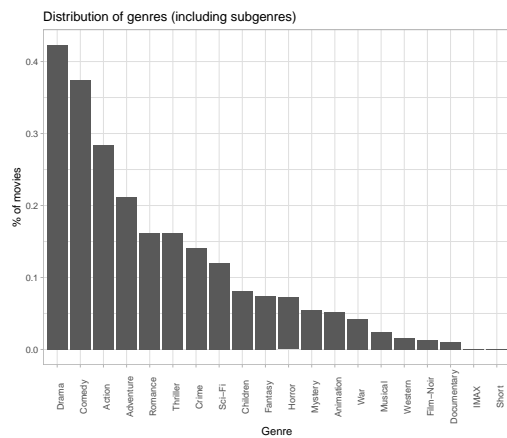


Most ratings are given as full number ratings (4, 3 or 5). Also the half-point ratings distribution follows a similar distribution as the full number ratings.

Now the distribution of the genres.



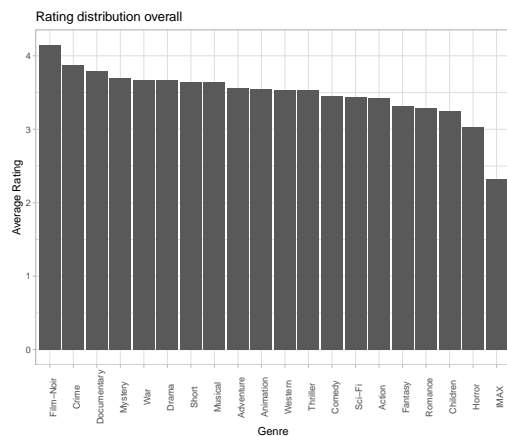
Most genres listed are Action, Comedy and Drama. What about when taking the sub genres into account.



When second and third listed genres are taken into account, the distribution is much finer. Top genres are still Drama, Comedy and Action but positions changed slightly. Drama is therefore an often used side genre, e.g. in combination with Romance, Thriller or others.

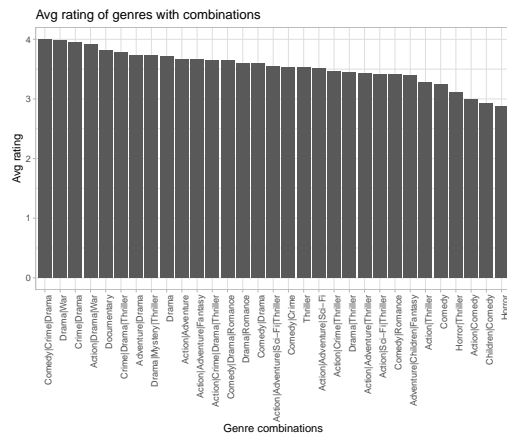
Genres

Lets plot the average rating against the genres.



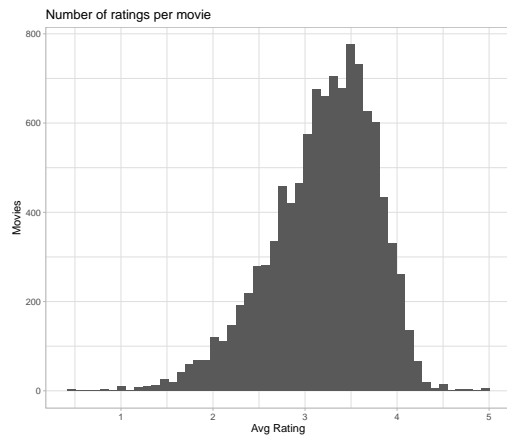
Ratings of genres show “intellectual” movie genres (e.g. Film-Noir, Crime, Drama..) are rated higher than movie genres associated with entertainment (e.g. Action, Fantasy, Horror)

And for average ratings of genre combinations with more than 50000 ratings.



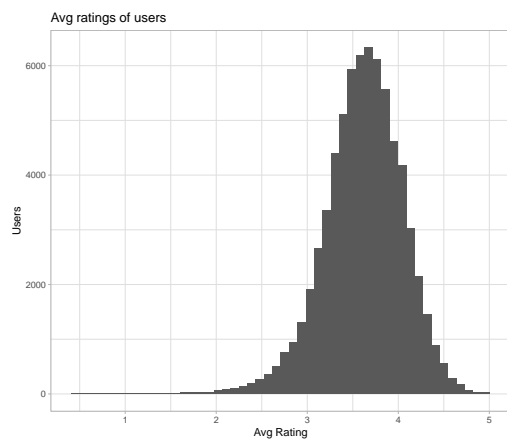
When taken genre combinations into account, a similar picture is painted, but some entertainment genre combinations, notably a Action combination, make it higher up the list (e.g. Action|Drama|War, Action|Adventure)

Average rating distribution of movies.



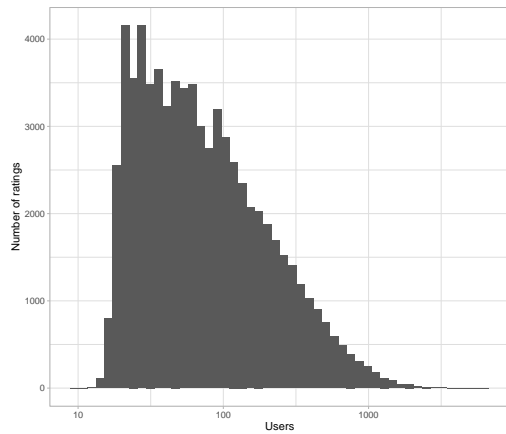
Only very few movies get an average rating above about 4.2. Most average ratings are between 2.5 and 4.

Average rating distribution of users



Most users rate movies between 3.5 and 4.2. This is higher than we've seen before on the average rating distribution of movies.

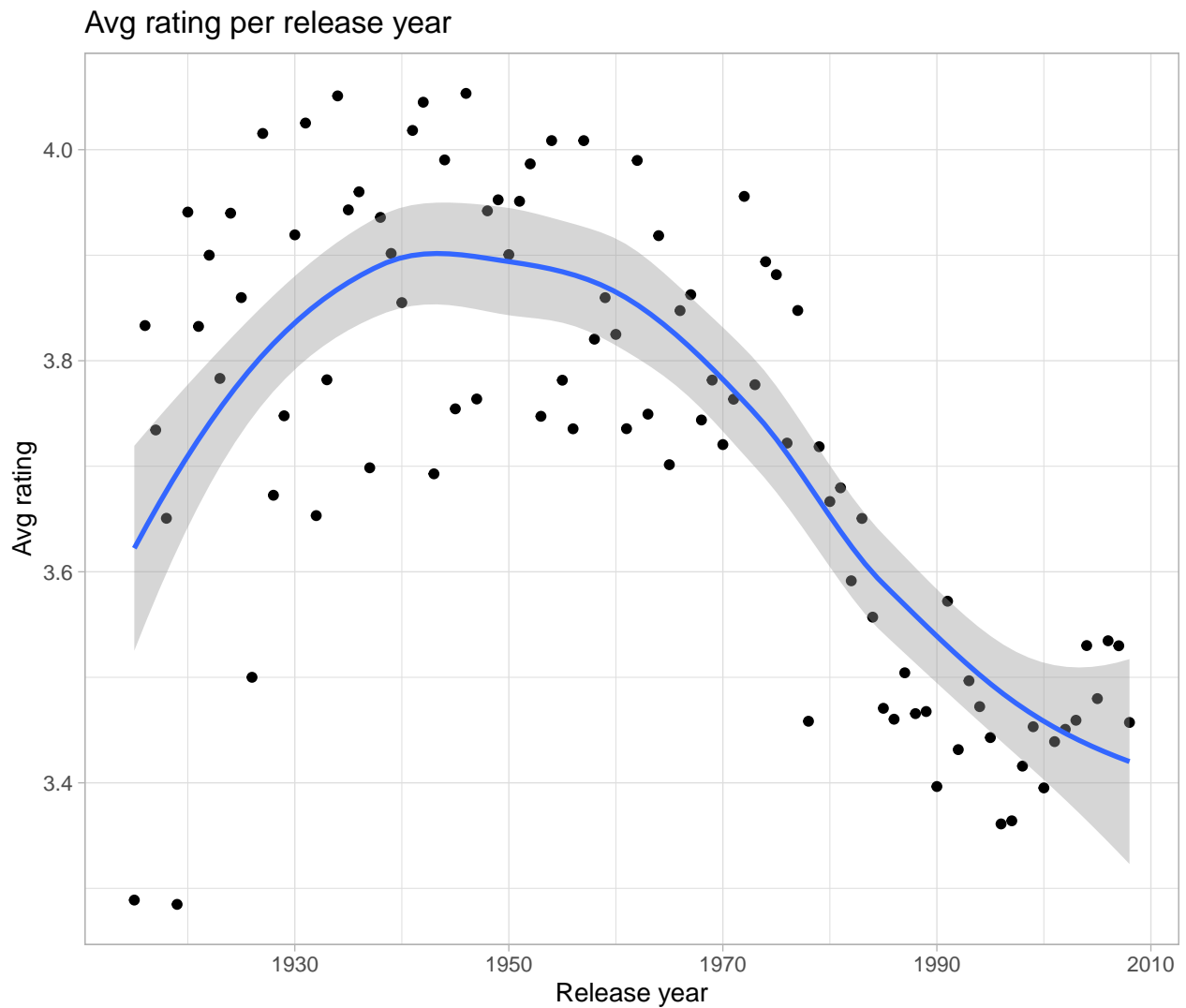
Number of ratings per user



Most users rate between a few and about 100 movies.

Average rating per release year with trendline

```
## `geom_smooth()` using method = 'loess' and formula = 'y ~ x'
```

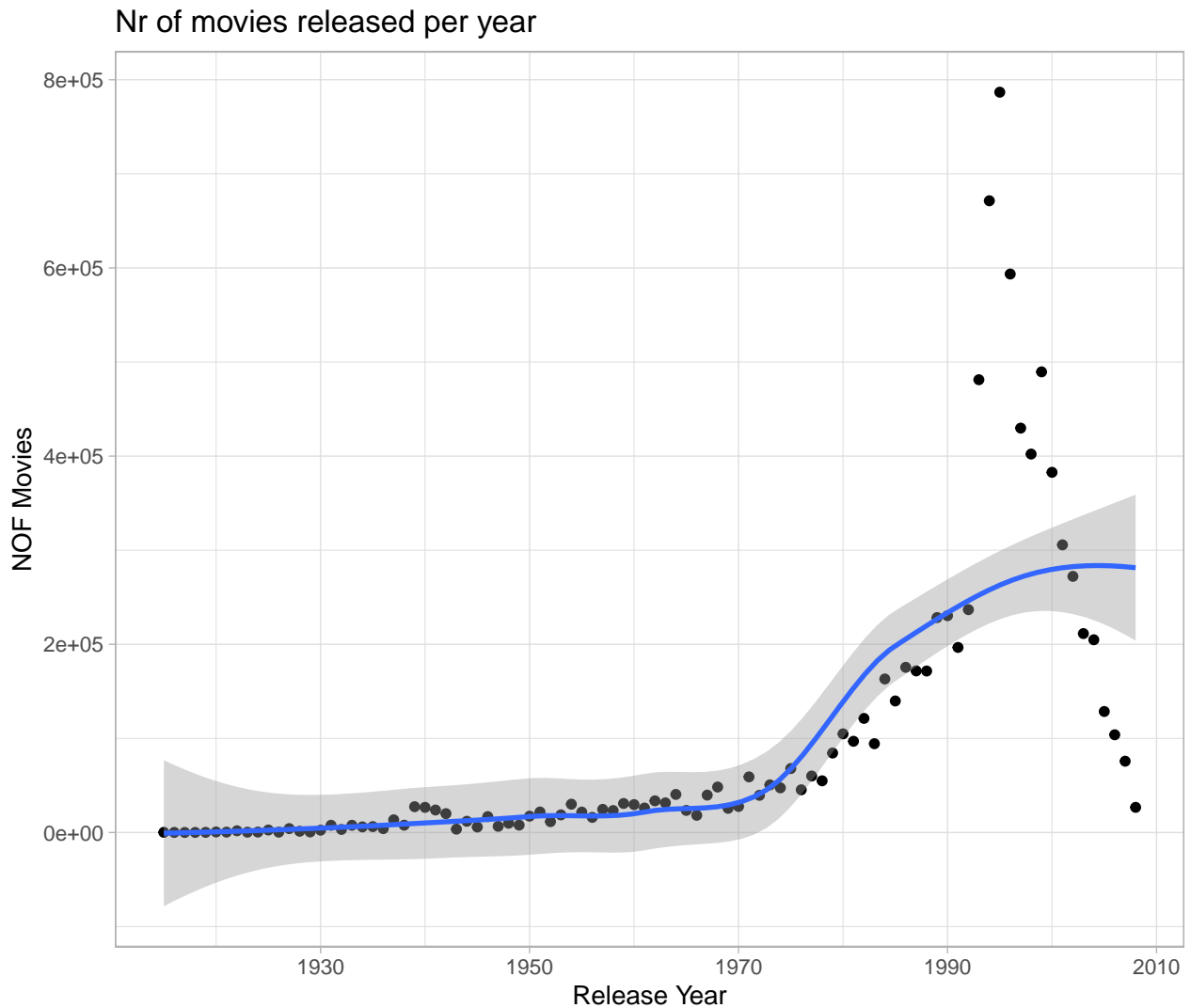


Movies released before 1970 are generally rated higher than movies released in the last two decades. This is a

know effect (only good stuff survive the test of time).

Number of released movies per year

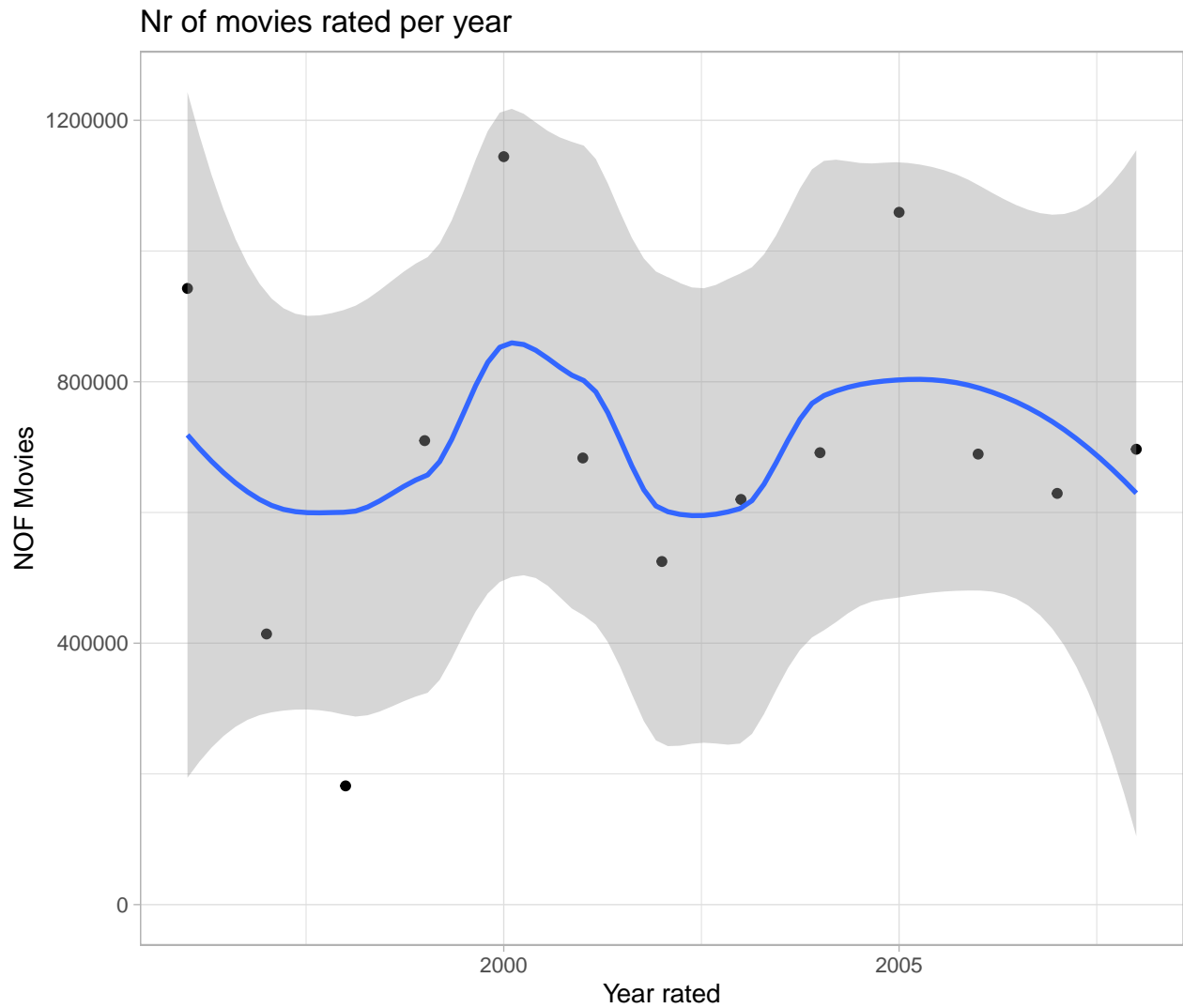
```
## `geom_smooth()` using method = 'loess' and formula = 'y ~ x'
```



Movie releases were relative stable before the 1970, then picked up and almost exploded 1990 and the following years, probably due to advancements in technology, production and market (e.g. movie theaters, movie rentals, tv...). Since before 2000 the number of movies released collapsed as quickly as its explosion a decade earlier.

Number of ratings per year. Only include the years from 1996 to 2008, data before and after is 0.

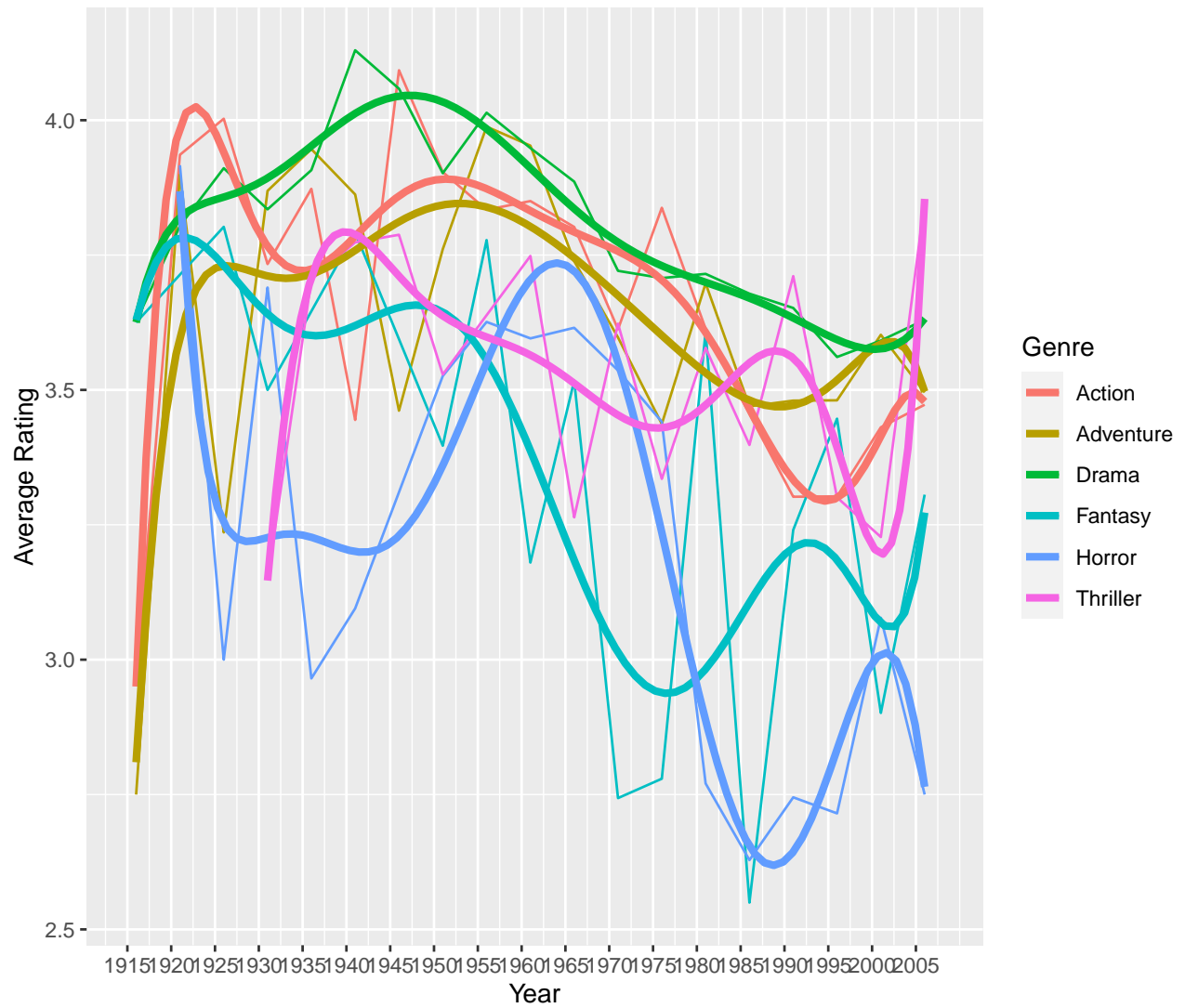
```
## `geom_smooth()` using method = 'loess' and formula = 'y ~ x'
```



Movie ratings per year are stable with some outliers.

Get the average rating of movie per genre but in 5 year bins.

```
## `summarise()` has grouped output by 'main_genre'. You can override using the  
## ``.groups` argument.
```

Some genres have pretty consistent average ratings over the years, others like e.g. Horror or Fantasy fluctuate a lot more.

Model

Prepare the training and test datasets and create a table for the RMSE results as we develop the model. I learned in the mean+movie approach, that all movieIds must be present in the training dataset, otherwise the test set will have movieId's where there was no training data on it. I guessed this will also be the case for userId and the main genre.

Guessing

Try stupid approach by guessing a rating from 0 to 5.

```
# create list of guessed ratings and make "guesses" for the size of the test set
guess_model <- sample(c(0, 0.5, 1, 1.5, 2, 2.5, 3, 3.5, 4, 4.5, 5), length(test_set$rating), replace=TRUE)
```

The resulting 2.1557864 is still far of the < 0.865, no surprise.

Model	RMSE
Guessing	2.155786

Mean

Get the average of all ratings on the training set and use this to predict a movie.

```
# Average mean of every movie in the training set
avg_model <- mean(train_set$rating)
```

The average of every movie used for predicting a rating results in 1.0598665. Closer but still far off.

Model	RMSE
Guessing	2.155786
Avg Model	1.059867

Genre bias

Evaluate the genre bias, the deviation from the movie average for each genre.

```
movie_avg <- mean(train_set$rating)

# bias of genres. deviation from the average for each genre. e.g. Film-Noir is 0.638
# above average, Horror 0.48 below average.
genre_bias <- train_set %>%
  group_by(main_genre) %>%
  summarise(deviation_genre = mean(rating - movie_avg))

# combine genre bias with the test_set
mean_genre_model <- test_set %>%
  inner_join(genre_bias, by="main_genre")
```

The RMSE is a bit better at 1.0480846 than just take the average like before.

Model	RMSE
Guessing	2.155786
Avg Model	1.059867
Genre Model	1.048085

Movie bias

The movie_bias is the difference of the avg movie rating to the mean rating.

```

movie_bias <- train_set %>%
  group_by(movieId) %>%
  summarise(deviation_movie = mean(rating - movie_avg))

# on the test set add the movie avg (3.512) with the difference the movie had
# to the avg in the training set and pull that column as a vector
mean_movie_model <- test_set %>%
  inner_join(movie_bias, by="movieId")

```

With and RMSE of 0.9431683 we are now in the sub 1 category.

Model	RMSE
Guessing	2.1557864
Avg Model	1.0598665
Genre Model	1.0480846
Movie Model	0.9431683

User and Movie bias

Let's add the average rating of a user into the mix.

```

# user_bias is the difference of the avg user rating to the mean rating
user_bias <- train_set %>%
  group_by(userId) %>%
  summarise(deviation_user = mean(rating - movie_avg))

mean_movie_user_model <- test_set %>%
  inner_join(movie_bias, by="movieId") %>%
  inner_join(user_bias, by="userId")

```

The user and movie gets us below 0.9. But 0.8854123 is still not near the desired < 0.865 .

Model	RMSE
Guessing	2.1557864
Avg Model	1.0598665
Genre Model	1.0480846
Movie Model	0.9431683
Movie + User Model	0.8854123

User, Movie and Genre bias

Now combine user, movie and genre together in a single model.

```

mean_movie_user_genre_model <- test_set %>%
  inner_join(movie_bias, by="movieId") %>%
  inner_join(user_bias, by="userId") %>%
  inner_join(genre_bias, by="main_genre")

```

This resulted in 0.9026264, which is worse than only using the movie and user. Maybe some tuning will fix it.

Model	RMSE
Guessing	2.1557864
Avg Model	1.0598665
Genre Model	1.0480846
Movie Model	0.9431683
Movie + User Model	0.8854123
Movie + User + Genre Model	0.9026264

Tuning

After searching for each tuning parameter individually these are the final tuning parameters:

- genre_bias: 0.0055
- movie_bias: 0.8944
- user_bias: 0.798

```
tuning_param <- seq(0, 1, 0.1)

rmse_seq <- sapply(tuning_param, function(t) {
  avg <- mean(train_set$rating)

  genre_bias <- train_set %>%
    group_by(main_genre) %>%
    summarise(deviation_genre = 0.0055 * sum(rating - movie_avg)/n())

  movie_bias <- train_set %>%
    group_by(movieId) %>%
    summarise(deviation_movie = 0.8944 * sum(rating - movie_avg)/n())

  user_bias <- train_set %>%
    group_by(userId) %>%
    summarise(deviation_user = 0.798 * sum(rating - movie_avg)/n())

  #release_year_1_bias <- train_set %>%
  #  group_by(releaseyear) %>%
  #  ifelse(releaseyear < 1930)
  #  summarise(deviation_releaseyear = t * sum(rating - movie_avg)/n())
  #
  #release_year_2_bias <- train_set %>%
  #  group_by(releaseyear) %>%
  #  summarise(deviation_releaseyear = t * sum(rating - movie_avg)/n())
  #
  #release_year_3_bias <- train_set %>%
  #  group_by(releaseyear) %>%
  #  summarise(deviation_releaseyear = t * sum(rating - movie_avg)/n())
  #
  #release_year_4_bias <- train_set %>%
  #  group_by(releaseyear) %>%
  #  summarise(deviation_releaseyear = t * sum(rating - movie_avg)/n())

  model <- test_set %>%
    inner_join(genre_bias, by="main_genre") %>%
    inner_join(movie_bias, by="movieId") %>%
    inner_join(user_bias, by="userId") # %>%
    #inner_join(release_year_bias, by="releaseyear")

  model$predicted_rating <- model$deviation_genre +
    model$deviation_user +
    model$deviation_movie +
    #model$deviation_releaseyear +
    movie_avg

  return(RMSE(test_set$rating, model$predicted_rating))
}
```

```
})
```

```
#qplot(tuning_param, rmse_seq, geom="line")
#param <- tuning_param[which.min(rmse_seq)]
#print(param)
```

first tried with only the genre_bias tuning, all other were $1 * \text{sum}(\text{rating} - \text{movie_avg})/n()$. This resulted in the best tuning parameter around 0. So the genre_bias effect is reduced having almost 0 effect. So in the end I removed it (commented it out for the history sake). -> genre_bias tune parameter: 0.0055

For the movie_bias testing the sequence seq(0, 1, 0.1) showed a minimum at about 0.9 -> movie_bias tune parameter: 0.8944

For the user_bias testing the sequence seq(0.5, 1.0, 0.1) showed a minimum at about 0.8 -> user_bias tune parameter: 0.798

Results

RMSE

Lets test the final model with its tuning in place against the verification set.

```
final_model_prediction <- function() {
  avg <- mean(train_set$rating)

  #genre_bias <- train_set %>%
  #  group_by(main_genre) %>%
  #  summarise(deviation_genre = 0.0055 * sum(rating - movie_avg)/n())

  movie_bias <- train_set %>%
    group_by(movieId) %>%
    summarise(deviation_movie = 0.8944 * sum(rating - movie_avg)/n())

  user_bias <- train_set %>%
    group_by(userId) %>%
    summarise(deviation_user = 0.798 * sum(rating - movie_avg)/n())

  model <- final_holdout_test %>%
    #inner_join(genre_bias, by="main_genre") %>%
    inner_join(movie_bias, by="movieId") %>%
    inner_join(user_bias, by="userId")

  model$predicted_rating <- #model$deviation_genre +
    model$deviation_user +
    model$deviation_movie +
    movie_avg
  return(model$predicted_rating)
}
```

This are all the achieved model RMSE:

Conclusion

Even with all this tuning, lower than 0.8786 is not possible with this approach.

- include year of rate or year diff from release to rating?

Model	RMSE
Guessing	2.1557864
Avg Model	1.0598665
Genre Model	1.0480846
Movie Model	0.9431683
Movie + User Model	0.8854123
Movie + User + Genre Model	0.9026264
Final Model Verification	0.8798800

- decade instead of release year? 5 year bins?
- side genres?
- remove outliers?

Future Improvements

Further information about the users could improve accuracy, e.g. shopping preferences, music taste, background information like education level. But there privacy concern about the usage of user personal data has to be considered. Training with larger dataset would be beneficial but would require more capable systems (e.g. with GPU).

Resources

[1] Rafael Irizarry. 2018. Introduction to Data Science.<https://rafalab.dfci.harvard.edu/dsbook/>

System

Hardware

All above computations are done with an Intel(R) Core(TM) i7-7500U CPU @ 2.70GHz CPU with 4 and 7.49 GB of RAM.

Software

This report is compiled using R markdown with RStudio.

```
sessionInfo()
```

```
## R version 4.2.3 (2023-03-15)
## Platform: x86_64-pc-linux-gnu (64-bit)
## Running under: Arch Linux
##
## Matrix products: default
## BLAS:   /usr/lib/libblas.so.3.11.0
## LAPACK: /usr/lib/liblapack.so.3.11.0
##
## Random number generation:
##  RNG:      Mersenne-Twister
##  Normal:   Inversion
##  Sample:   Rounding
##
## locale:
##  [1] LC_CTYPE=en_GB.UTF-8      LC_NUMERIC=C
##  [3] LC_TIME=en_GB.UTF-8      LC_COLLATE=en_GB.UTF-8
##  [5] LC_MONETARY=en_GB.UTF-8  LC_MESSAGES=en_GB.UTF-8
##  [7] LC_PAPER=en_GB.UTF-8     LC_NAME=C
##  [9] LC_ADDRESS=C             LC_TELEPHONE=C
## [11] LC_MEASUREMENT=en_GB.UTF-8 LC_IDENTIFICATION=C
##
## attached base packages:
## [1] parallel  stats      graphics  grDevices  utils      datasets  methods
## [8] base
##
## other attached packages:
##  [1] caret_6.0-93      lattice_0.20-45    forcats_0.5.2      dplyr_1.0.10
##  [5] purrr_0.3.5       readr_2.1.3        tidyr_1.2.1        tibble_3.2.0
##  [9] tidyverse_1.3.2    kableExtra_1.3.4   stringr_1.4.1      scales_1.2.1
## [13] lubridate_1.9.0    timechange_0.1.1   rmarkdown_2.18     benchmarkme_1.0.8
## [17] devtools_2.4.5     usethis_2.1.6      ggplot2_3.4.0      pacman_0.5.1
##
## loaded via a namespace (and not attached):
##  [1] googledrive_2.0.0    colorspace_2.0-3    ellipsis_0.3.2
##  [4] class_7.3-21         fs_1.5.2            rstudioapi_0.14
##  [7] farver_2.1.1         listenv_0.8.0       remotes_2.4.2
## [10] bit64_4.0.5          prodlim_2019.11.13  fansi_1.0.4
## [13] xml2_1.3.3           codetools_0.2-19    splines_4.2.3
```

## [16]	doParallel_1.0.17	cachem_1.0.6	knitr_1.41
## [19]	pkgload_1.3.2	jsonlite_1.8.4	pROC_1.18.0
## [22]	broom_1.0.1	dbplyr_2.2.1	shiny_1.7.4
## [25]	compiler_4.2.3	httr_1.4.5	backports_1.4.1
## [28]	assertthat_0.2.1	Matrix_1.5-3	fastmap_1.1.0
## [31]	gargle_1.2.1	cli_3.6.0	later_1.3.0
## [34]	htmltools_0.5.4	prettyunits_1.1.1	tools_4.2.3
## [37]	gtable_0.3.1	glue_1.6.2	reshape2_1.4.4
## [40]	Rcpp_1.0.9	cellranger_1.1.0	vctrs_0.5.2
## [43]	svglite_2.1.1	nlme_3.1-162	iterators_1.0.14
## [46]	timeDate_4021.106	xfun_0.35	gower_1.0.0
## [49]	globals_0.16.2	ps_1.7.2	rvest_1.0.3
## [52]	mime_0.12	miniUI_0.1.1.1	lifecycle_1.0.3
## [55]	googlesheets4_1.0.1	future_1.29.0	MASS_7.3-58.2
## [58]	ipred_0.9-13	vroom_1.6.0	hms_1.1.2
## [61]	promises_1.2.0.1	yaml_2.3.6	memoise_2.0.1
## [64]	rpart_4.1.19	stringi_1.7.8	highr_0.9
## [67]	foreach_1.5.2	hardhat_1.2.0	pkgbuild_1.4.0
## [70]	lava_1.7.0	benchmarkmeData_1.0.4	rlang_1.0.6
## [73]	pkgconfig_2.0.3	systemfonts_1.0.4	archive_1.1.5
## [76]	evaluate_0.18	labeling_0.4.2	recipes_1.0.3
## [79]	htmlwidgets_1.6.1	bit_4.0.5	processx_3.8.0
## [82]	tidyselect_1.2.0	parallelly_1.32.1	plyr_1.8.8
## [85]	magrittr_2.0.3	R6_2.5.1	generics_0.1.3
## [88]	profvis_0.3.7	DBI_1.1.3	mgcv_1.8-42
## [91]	pillar_1.8.1	haven_2.5.1	withr_2.5.0
## [94]	survival_3.5-3	nnet_7.3-18	future.apply_1.10.0
## [97]	modelr_0.1.10	crayon_1.5.2	utf8_1.2.3
## [100]	tzdb_0.3.0	urlchecker_1.0.1	grid_4.2.3
## [103]	readxl_1.4.1	data.table_1.14.6	callr_3.7.3
## [106]	ModelMetrics_1.2.2.2	reprex_2.0.2	digest_0.6.30
## [109]	webshot_0.5.4	xtable_1.8-4	httpuv_1.6.8
## [112]	stats4_4.2.3	munsell_0.5.0	viridisLite_0.4.1
## [115]	sessioninfo_1.2.2		

Harper, F. Maxwell, and Joseph A. Konstan. 2015. “The MovieLens Datasets.” *ACM Transactions on Interactive Intelligent Systems* 5 (4): 1–19. <https://doi.org/10.1145/2827872>.